



# Scene classification using local and global features with collaborative representation fusion



Jinyi Zou<sup>a</sup>, Wei Li<sup>a,\*</sup>, Chen Chen<sup>b</sup>, Qian Du<sup>c</sup>

<sup>a</sup> College of Information Science and Technology, Beijing University of Chemical Technology, Beijing 100029, China

<sup>b</sup> Department of Electrical Engineering, University of Texas at Dallas, Richardson, TX 75080, USA

<sup>c</sup> Department of Electrical and Computer Engineering, Mississippi State University, MS 39762, USA

## ARTICLE INFO

### Article history:

Received 2 November 2015

Revised 4 February 2016

Accepted 8 February 2016

Available online 13 February 2016

### Keywords:

Scene classification

Locality-constrained linear coding

Spatial pyramid matching

Collaborative representation-based classification

## ABSTRACT

This paper presents an effective scene classification approach based on collaborative representation fusion of local and global spatial features. First, a visual word codebook is constructed by partitioning an image into dense regions, followed by the typical  $k$ -means clustering. A locality-constrained linear coding is employed on dense regions via the visual codebook, and a spatial pyramid matching strategy is then used to combine local features of the entire image. For global feature extraction, the method called multiscale completed local binary patterns (MS-CLBP) is applied to both the original gray scale image and its Gabor feature images. Finally, kernel collaborative representation-based classification (KCRC) is employed on the extracted local and global features, and class label of the testing image is assigned according to the minimal approximation residual after fusion. The proposed method is evaluated by using four commonly-used datasets including two remote sensing images datasets, an indoor and outdoor scenes dataset, and a sports action dataset. Experimental results demonstrate that the proposed method significantly outperforms the state-of-the-art methods.

© 2016 Elsevier Inc. All rights reserved.

## 1. Introduction

In the last decade, scene classification has drawn increasing attention both in academia and industry [14,37,45,46,57,63]. The task is to automatically classify an image by feature extraction and label assignment. Although great effort in extracting features (e.g., hash codes [10,15,17], manifold structures [24,56,58], etc) has been made, it is still a challenging task due to many factors to be considered such as variations in spatial position, illumination, and scale.

In the early days, scene classification methods mainly concentrated on modeling [25] and using global spatial features such as color and texture histograms [41]. The global features often have simple implementation and low computational cost but offer limited performance. In [11,36,48,50,59,61], the popular bag-of-visual-words (BoVW) model was adopted, which represented an image with an orderless collection of local features. In this model, an image can be treated as a document, similar to “words”. The image is usually partitioned into patches and represented by a codebook. To this end, it follows three steps: (i) feature detection (commonly, the key point detection [30,31] is used.), (ii) feature description based on key points [9,27], and (iii) codebook generation. However, the BoVW model ignores the spatial layout of features.

\* Corresponding author. Tel.: +86 10 64413467, +86 18146529853; fax: +86 10 64434726.

E-mail address: [liwei089@ieee.org](mailto:liwei089@ieee.org), [leewei36@gmail.com](mailto:leewei36@gmail.com) (W. Li).

To overcome this issue, spatial pyramid matching (SPM) was developed based on approximate global geometric correspondence in [16]. This approach partitions an image into increasingly fine sub-regions and computes histograms of local features in each sub-region. It is a simple extension of an orderless bag-of-features image representation and overcomes the shortcoming of the BoVW model. Based on this SPM framework, many extensions have been proposed. In [52], an improved SPM method was developed by generalizing vector quantization (VQ) to sparse coding, followed by multiscale spatial max pooling. A linear SPM kernel based on Lowes scale invariant feature transform (SIFT) [9,27] sparse codes was employed, providing excellent performance on several scene datasets. However, this method exhibits high computation complexity and is extremely time-consuming. In [49], a simple yet effective coding scheme called locality-constrained linear coding (LLC) was proposed to replace the VQ coding. LLC utilizes the locality constraint to project each descriptor onto its local-coordinate system and applies max pooling to the projected coordinates to generate the final representation. Compared with the sparse coding in [3,52], LLC not only guarantees sparsity but also solves the representation problem as a constrained least squares fitting problem, which takes both accuracy and efficiency into consideration.

With the success of the BoVW and SPM models in image scene classification, more algorithms have been developed based on these frameworks. For example, a multi-resolution representation was incorporated into the BoVW model [69,70], which constructed images with multiple resolutions and extracted local features from all the images with dense regions, utilized the  $k$ -means clustering to generate visual codebook, and then represented each sub-region as a histogram of code-word occurrences by mapping the local features to the codebook. In [8], the pyramid histogram of multi-scale block local binary pattern (PH-MBLBP) descriptor was employed, which can encode micro-and macro-structures of images, and the PH-MBLBP descriptor was verified to be a powerful texture descriptor with low computational complexity. In [66], a concentric circle-based spatial-rotation-invariant representation strategy for describing spatial information of visual words and a concentric circle-structured multi-scale BoVW method using multiple features were proposed to enhance image rotation invariance, leading to excellent scene classification results. In [65], a scene image was transformed into multi-features by 2-D wavelet decompositions, and the resulting feature maps were then employed by the BoVW and SPM models; after that, all the features of different feature maps were stacked as inputs to a classic support vector machine (SVM) classifier.

Although the BoVW [54] and SPM, [6,55] models are popular in scene classification and have gained great success, one of disadvantages is that they cannot well capture global structures in an image scene. Compared with the BoVW and SPM models, some global feature based scene classification methods actually have achieved satisfying performance. In [4], a global feature called completed local binary pattern (CLBP) was employed and multi-scale resolution was adopted to generate multi-scale global features for land-use scene classification. In [5], a global Gabor-filtering-based CLBP was applied to generate multi-features. Some other global feature representation methods for scene classification can be found in [19,20,38,40,68]. It is known that the BoVW model with SPM can capture spatial information with ordered block partitions but it is sensitive to rotation variations. On the other hand, the global feature representation focuses on global texture, such as texture depth and global contrast, but ignores local details or small objects.

Since each type of features (i.e., local features such as BoVW, and global features such as CLBP) has its own advantages and limitations, we propose to fuse the global and local features together to characterize both local fine details and global structures of scene images. In traditional methods, the fusion strategy mainly includes feature and decision level fusion [67]. Nevertheless, it is difficult to combine different features in feature level because different features may not be compatible. Moreover, it may impose a challenge of high feature dimensionality in feature level fusion. In decision level fusion, voting may lead to a rough result. Thus, in this work, both feature and decision level-fusion methods are considered to mitigate their individual shortcomings in the proposed framework. We first use the BoVW and SPM to generate local features and employ the multi-scale CLBP (MS-CLBP) to extract global features. We then employ the feature representation method (e.g., kernel collaborative representation-based classification, KCRC [53]) which is different from a conventional training-testing classifier (e.g., SVM or extreme learning machine (ELM) [22]). The features in training images can be presented as a dictionary, then testing images are reconstructed by the training dictionary. After obtaining the representation residuals from using two types of features, the sum of weighted residuals is calculated and the label is assigned according to the minimal residual class. Experimental results on four benchmark datasets demonstrate that our approach gains a remarkable classification improvement over the state-of-the-art methods.

The main contributions made in this paper can be summarized as follows. (1) To the best of our knowledge, it is the first time that local and global features are employed together for scene image classification. The complementary nature of these two types of features can effectively mitigate the shortcomings of local feature representation based methods (e.g., BoVW and SPM) and global feature representation based methods (e.g., MS-CLBP). (2) A weighted sum of approximation residuals of local and global features with collaborative representation fusion are proposed, which can overcome the difficulties in feature or decision level fusion.

The rest of the paper is organized as follows. Section 2 describes the proposed approaches, including local and global feature extraction and the collaborative representation fusion strategy. Section 3 introduces four experimental datasets. Section 4 represents the experimental results and provides some discussion. Section 5 draws the conclusions.

## 2. Proposed approach

The flow chart of the proposed method is illustrated in Fig. 1. First, the local and global feature dictionaries are prepared in advance. Then, for each testing image, the local and global features are extracted and collaboratively represented by the

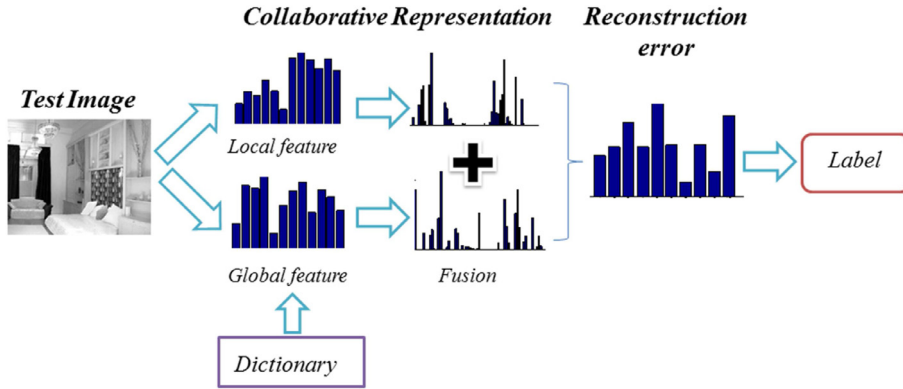


Fig. 1. Overview of the local and global features with collaborative representation fusion based scene classification.

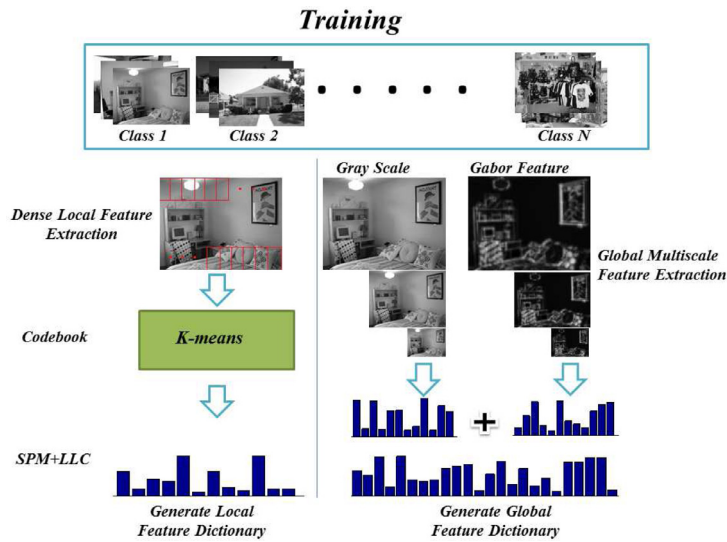


Fig. 2. Overview of the procedure of generating the local and global training dictionaries.

training dictionary. The dictionary is constructed by the specific classes and the training samples of each class can be viewed as a sub-dictionary. Through collaborative representation, local and global features are coded by the sub-dictionaries. The coefficients of each sub-dictionary present the contribution of each class to the testing image. The reconstruction errors of each class measure the similarity of the testing class and the training classes. Both the local and global features are used for representation and two reconstruction residuals are fused. The label of the class with the minimal reconstruction residual is then assigned to the testing image. In this way, both the local and global information can be utilized even when some images mainly exhibit global information (e.g., beach scene and desert scene), or local information (e.g., indoor object scene). The proposed method can deal with scene images with different levels of details by considering these two types of features adaptively.

### 2.1. Feature extraction and dictionary generation

In recent years, visual descriptors have shown their outstanding performance in image classification. For instance, SIFT [9,27] exhibits its powerful description capability in object recognition, detection, and image matching, etc [28,42]. In [6,16,18], it has been demonstrated that the dense feature extraction outperforms better than the interest point based feature extraction. This is because the interest points detected by Harris corner detector or difference of Gaussian pyramid are always located at areas with great changes; however, smooth areas could also contain important scene information.

In this work, we employ the dense SIFT descriptor by partitioning an image into small patches. We refer to [49] for local feature generation and adopt the MS-CLBP [4] for global features. The detailed process is presented in Fig. 2. Specifically, images from each class are randomly chosen to perform dense local feature extraction. A regular grid is used to partition an entire image into patches. For each patch, the patch center is viewed as a key point, on which the SIFT feature is generated. After that, *k*-means clustering is employed to generate the codebook which presents the visual words of the BoVW. Then,

vector quantization is used to encode the patches using the codebook and calculate the histogram of frequency of each visual word. The SPM and LLC are further employed to calculate the local features due to that the SPM is able to avoid the orderless bag-of-features representation which enhances the spatial order structure and leads to better representation of local features. In addition, the LLC is used for reconstruction instead of direct vector quantization since it can reconstruct the feature with high accuracy and low computational cost.

Let  $\mathbf{X}$  be a set of  $d$ -dimensional local descriptors extracted from an image with  $n$  entries, i.e.,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in R^{d \times n}$ , and let the codebook with  $m$  entries be denoted as  $\mathbf{D} = [\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_m] \in R^{d \times m}$ . The vector quantization solves the following constrained least squares fitting problem,

$$\arg \min_{\alpha} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{D}\alpha_i\|^2 \quad \text{subject to} \quad \|\alpha_i\|_0 = 1, \|\alpha_i\|_1 = 1, \alpha_i \geq 0 \quad (1)$$

where  $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_n]$  is a coding coefficient vector for  $\mathbf{X}$  and  $\alpha_i \in R^m \times 1$ . For each patch descriptor, the vector quantization only finds the nearest atom from the codebook. It may cause large reconstruction errors and lead to class mis-assignment. A better coding strategy is the sparse representation [26] or the collaborative representation [53]. Here, we do not consider sparse coding but the LLC because the former requires large computational effort. The objective function of the LLC can be expressed as

$$\arg \min_{\alpha} \|\mathbf{x}_i - \mathbf{D}\alpha_i\|^2 + \lambda \|\mathbf{l}_i \odot \alpha_i\|^2 \quad \text{subject to} \quad \mathbf{1}^T \alpha_i = 1 \quad (2)$$

where  $\odot$  denotes the element-wise multiplication, and  $\mathbf{l}_i$  denotes the  $i$ th locality adaptor that assigns weights to each atom according to its similarity to the input descriptor. Define  $\mathbf{l}_i = f(\mathbf{x}_i, \mathbf{D})$ , where  $f(\cdot)$  measures the similarity between the input descriptor and the atoms of codebook. To achieve sparsity of the LLC, only the most similar atoms are allowed to have nonzero coefficients. For each training image, local features can be obtained. Thus, the training local feature dictionary consists of the local features from all the training images.

On the other hand, a global feature dictionary is generated based on CLBP [4,13] which is a completed modeling of the typical local binary pattern (LBP) [33]. LBP is an effective measure of spatial structure information and has been widely employed in face recognition and texture classification [1,35]. Given a center pixel  $\mathbf{y}_c$  and its neighboring pixels equally distributed on a circle of radius  $r$ . If the coordinates of  $\mathbf{y}_c$  are  $(0, 0)$  and the coordinates of  $m$  neighbors  $\{\mathbf{y}_i\}_{i=0}^{m-1}$  are  $(-r \sin(2\pi i/m), r \cos(2\pi i/m))$ . The LBP is computed by thresholding the neighbors with the center pixel to generate an  $m$ -bit binary code,

$$LBP_{m,r}(\mathbf{y}_c) = \sum_{i=0}^{m-1} s(\mathbf{y}_i - \mathbf{y}_c) \times 2^i. \quad (3)$$

If  $\mathbf{y}_i - \mathbf{y}_c \geq 0$ ,  $s(\mathbf{y}_i - \mathbf{y}_c) = 1$ ; otherwise,  $s(\mathbf{y}_i - \mathbf{y}_c) = 0$ . However, the LBP only uses the sign information and ignores the magnitude information. CLBP is designed to combine both the sign and magnitude information for comprehensive modeling. CLBP-Sign (CLBP\_S) is equivalent to the traditional LBP and the CLBP-Magnitude (CLBP\_M) is defined as,

$$CLBP\_M_{m,r} = \sum_{i=0}^{m-1} p(|\mathbf{y}_i - \mathbf{y}_c|, \gamma) \times 2^i, \quad \text{and} \quad p(\delta, \gamma) = \begin{cases} 1, & \delta \geq \gamma \\ 0, & \delta < \gamma \end{cases} \quad (4)$$

where  $\gamma$  is the mean value of  $|\mathbf{y}_i - \mathbf{y}_c|$  from the entire image. Then, an  $m$ -bit binary code can also be generated. The CLBP\_S is able to provide the spatial texture and the CLBP\_M is mainly to describe the depth of texture. Note that we do not use the CLBP-Center component in CLBP which encodes the pixel intensity for the purpose of reducing computational complexity. The syncretic CLBP can describe both the spatial and depth information by mapping the CLBP\_S and CLBP\_M into histograms. To further enhance the representation power, the multi-scale representation of CLBP is considered. Specifically, the original image is down-sampled using the bicubic interpolation to obtain multiple images with different sizes (or scales). Then, the CLBP operator is applied with the same parameters to these images at different scales. Finally, the CLBP features are stacked to form the global features. In addition, it has been verified that the Gabor features with LBP show excellent performance in face recognition [64]. Therefore, global feature extraction is performed (as shown in Fig. 2) on the Gabor feature images in addition to the original intensity images. As a result, MS-CLBP features from gray-level images and the corresponding Gabor feature images are combined as the final global features. The training global feature dictionary can thus be obtained.

## 2.2. Image classification by fusing reconstruction residuals

After generating the local and global feature dictionaries, KCRC [53] is employed to represent a testing image using these features. Furthermore, a weighted fusion strategy is considered for class label determination based on reconstruction errors (residuals) as shown in Fig. 3, where the image is first represented by training feature dictionary and then reconstructed by each sub-dictionary. The middle histogram represents the residual of reconstruction error of each class according to each feature. Fig. 3 further illustrates that the minimum residual and the second smallest value may be very close, introducing classification error. However, after the local and global feature reconstruction errors are fused, the discrimination power is enhanced.

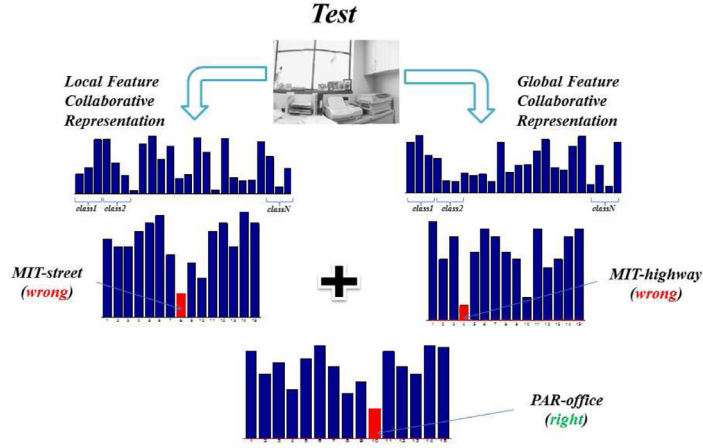


Fig. 3. Illustration of representation residual fusion.

For a testing image, local features are encoded via the training local feature dictionary by KCRC to obtain coding coefficients. The same operation is applied to global features. Before introducing the KCRC, we first briefly overview the original collaborative representation-based classification (CRC) [23,62]. The KCRC can be viewed as a kernel version of CRC, which codes a testing sample by a linear combination of all the training samples. Let  $\mathbf{y}$  be the feature of a testing image and the training dictionary be  $\mathbf{X} = [\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^C]$  with  $C$  classes, where  $\mathbf{X}^l$  is a sub-dictionary for the  $l$ th class. The objective function of CRC is,

$$\mathbf{w} = \arg \min \{ \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \}, \quad (5)$$

where  $\mathbf{w}$  is the weight vector and  $\lambda$  is the regularization parameter. Eq. (5) can be directly solved as

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}. \quad (6)$$

After obtaining  $\hat{\mathbf{w}}$ ,  $\mathbf{X}$  and  $\hat{\mathbf{w}}$  are separated into  $C$  class-specific segments, i.e.,  $\{\mathbf{X}^l\}_{l=1}^C$  and  $\{\hat{\mathbf{w}}^l\}_{l=1}^C$ . The class-specific residuals can be calculated as,

$$r^l(\mathbf{y}) = \|\mathbf{y} - \mathbf{X}^l \hat{\mathbf{w}}^l\|_2, \quad l = 1, 2, \dots, C. \quad (7)$$

The label of the testing image is then determined according to

$$\text{Class}(\mathbf{y}) = \arg \min_{l=1, \dots, C} r^l(\mathbf{y}). \quad (8)$$

The CRC is a linear method in nature and may not handle samples with complex nonlinear variations due to pose, expression, and illumination. Nevertheless, the kernel-based methods [32,51] can effectively discover nonlinear structures by mapping the samples into a high dimensional space.

In this work, the KCRC algorithm is adopted to solve the nonlinear reconstruction problem. In a kernel method, data are projected into a high-dimensional kernel-induced feature space by an implicit nonlinear mapping function  $\Phi$ , i.e.,  $\mathbf{y} \rightarrow \Phi(\mathbf{y}) \in \mathbb{R}^D \times 1$  ( $D \gg d$  is the dimension of kernel feature space). Usually, the explicit mapping  $\Phi$  is unknown. According to the kernel trick, the inner product of two vectors,  $\mathbf{x}$  and  $\mathbf{y}$ , in the feature space is equal to the output of a kernel function, i.e.,  $k(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x})^T \Phi(\mathbf{y})$ . Here, we employ the Gaussian radial basis function (RBF) kernel  $k(\mathbf{x}, \mathbf{y}) = \exp(-\sigma \|\mathbf{x} - \mathbf{y}\|_2^2)$  ( $\sigma > 0$  is a parameter of the RBF kernel).

In KCRC, the weight vector  $\mathbf{w}$  is estimated as

$$\mathbf{w} = \arg \min \|\Phi(\mathbf{y}) - \Phi \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2, \quad (9)$$

where  $\Phi = [\Phi(\mathbf{x}_1), \Phi(\mathbf{x}_2), \dots, \Phi(\mathbf{x}_n)] \in \mathbb{R}^{D \times n}$ , and the closed-form solution is

$$\hat{\mathbf{w}} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{k}(\cdot, \mathbf{y}), \quad (10)$$

where  $\mathbf{k}(\cdot, \mathbf{y}) = [k(\mathbf{x}_1, \mathbf{y}), k(\mathbf{x}_2, \mathbf{y}), \dots, k(\mathbf{x}_n, \mathbf{y})]^T \in \mathbb{R}^{n \times 1}$ , and  $\mathbf{K} = \Phi^T \Phi \in \mathbb{R}^{n \times n}$  is the Gram matrix with  $\mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$ .

After computing the weight vector  $\hat{\mathbf{w}}$  in the kernel space, the class-specific error residual can be calculated as

$$\begin{aligned} r^l(\mathbf{y}) &= \|\Phi^l \hat{\mathbf{w}}^l - \Phi(\mathbf{y})\|_2 \\ &= \sqrt{(\Phi(\mathbf{y}) - \Phi^l \hat{\mathbf{w}}^l)^T (\Phi(\mathbf{y}) - \Phi^l \hat{\mathbf{w}}^l)} \\ &= \sqrt{k(\mathbf{y}, \mathbf{y}) + (\hat{\mathbf{w}}^l)^T \mathbf{K}^l \hat{\mathbf{w}}^l - 2(\hat{\mathbf{w}}^l)^T \mathbf{k}^l(\cdot, \mathbf{y})}, \end{aligned} \quad (11)$$



Fig. 4. Example images from the 21-class land-use dataset.

where  $\Phi^l$  represents the kernel sub-dictionary for class  $l$  and  $\mathbf{K}^l$  is the Gram matrix of the samples in class  $l$ . Finally, class label of the testing image is determined using Eq. (8).

In the proposed local-global-fusion (LGF) strategy, each testing image is generated two types of features, i.e.,  $\mathbf{y}_L$  and  $\mathbf{y}_G$  representing local and global features, respectively. The KCRC is employed to calculate the approximation residuals of both local and global features denoted by  $r_L^l(\mathbf{y})$  and  $r_G^l(\mathbf{y})$ , and the local and global residuals are fused as,

$$r_{LGF}^l(\mathbf{y}) = \mu r_L^l(\mathbf{y}) + (1 - \mu)r_G^l(\mathbf{y}), \quad (12)$$

where  $\mu$  is the weighting parameter. Finally, class label of the testing image is determined using Eq. (8) based on  $r_{LGF}^l(\mathbf{y})$ . As illustrated in Fig. 3, individual decisions based on residuals from local or global features are wrong, but the proposed LGF method generates a correct result after residual fusion.

### 3. Experimental datasets

The first dataset is the well-known UC Merced Land Use dataset [54] consisting of 21 land-use classes, i.e., agricultural, airplane, baseball diamond, beach, buildings, chaparral, dense residential, forest, freeway, golf course, harbor, intersection, medium density residential, mobile home park, overpass, parking lot, river, runway, sparse residential, storage tanks, and tennis courts. Each class contains 100 images of size  $256 \times 256$  pixels. It is up to now the first public land-use scene image dataset with ground truth, which is created by Yang and Newsam<sup>1</sup>. This dataset was downloaded from United States Geological Survey (USGS) National Map. Sample images of each land-use class are illustrated in Fig. 4. This benchmark dataset has a large geographical scale. In order to facilitate a fair comparison, we follow the same experimental setup reported in [54]. We randomly choose 80 images from each land-use class for training and the remaining for testing. All the experiments are repeated 10 trials with different realizations of training and testing images.

The second dataset used in our experiments is the 19-class satellite scene dataset [4,44]. It consists of 19 classes of high-resolution satellite scenes including airport, beach, bridge, commercial, desert, farmland, football field, forest, industrial, meadow, mountain, park, parking, pond, port, railway station, residential, river and viaduct. There are 50 images of size  $600 \times 600$  pixels for each class. Sample images of each class are shown in Fig. 5. This benchmark dataset also focuses on images with a large geographical scale. We follow the same experimental setup in [4]. We randomly choose 30 images from each class for training and the remaining for testing. All the experiments are repeated 10 trials.

The third dataset includes the 15-scene categories [8]. The initial 8 classes were collected by Oliva and Torralba [34], and the 5 categories were added by Li and Perona [18], the rest categories were introduced by Lazebnik [16]. The 15 scene categories are: bedroom (216 images), CAL-suburb (241 images), industrial (311 images), kitchen (210 images), livingroom (289 images), MIT-cost (360 images), MIT-forest (328 images), MIT-highway (260 images), MIT-insidecity (308 images), MIT-mountain (374 images), MIT-opencountry (410 images), MITstreet (292 images), MIT-tallbuilding (356 images), PAR-office (215 images), and store (315 images). It contains a total of 4485 images of size  $300 \times 250$  pixels. This is the most widely-used dataset in literature so far. It is a challenging dataset since it contains a wide range of indoor and outdoor images. Both local and global scales images are included. Sample images of each class are shown in Fig. 6. Following the same experimental procedure [16], we randomly choose 100 images per class for training and the rest for testing.

The fourth dataset has sports event categories collected from the Internet [19]. The eight sports event categories are: bocce (137 images), croquet (236 images), polo (182 images), rowing (250 images), snowboarding (190 images), badminton

<sup>1</sup> <http://vision.ucmerced.edu/datasets>.

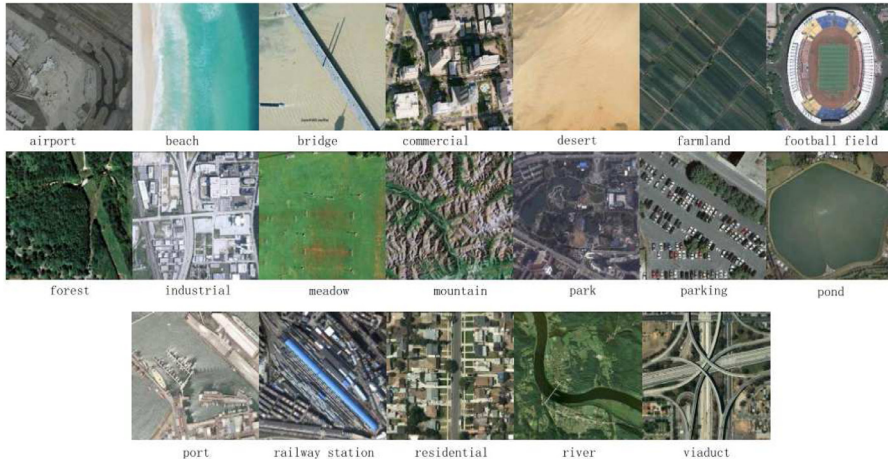


Fig. 5. Example images from the 19-class satellite scene dataset.

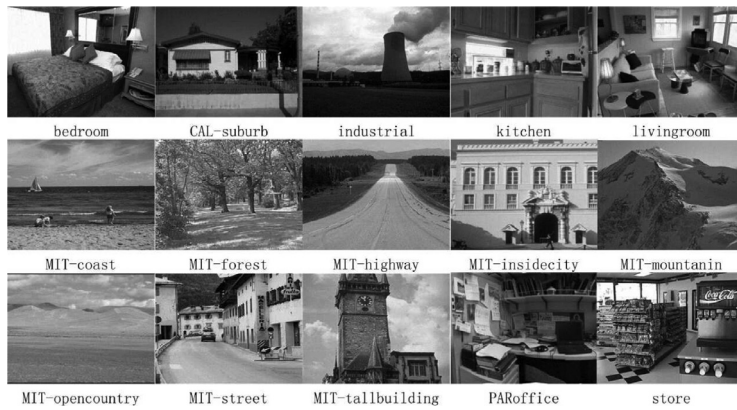


Fig. 6. Example images from the 15-scene categories dataset.

(200 images), sailing (190 images), and rock climbing (194 images). Sample images of each class are shown in Fig. 7. This event dataset is also very challenging mainly because: (1) the background of images is highly cluttered and diverse, (2) the object classes are also diverse, and (3) the same category sizes of instances from the same object are different. In addition, this dataset contains a large range of image scales. Following the same experimental setting in [19], 70 training samples are randomly selected from each class and the remaining ones for testing.

#### 4. Experimental analysis

For the UC-Merced land-cover scene dataset, the suggested parameters in [52] are used to generate the descriptor codebook. The dense SIFT descriptors are extracted for each  $16 \times 16$  patch on a grid of step size of 6 pixels. The codebook size is set to be 500. There are 1680 images for training (80 per class). Therefore, a total of 2,824,080 patches is generated. For  $k$ -means clustering, 30,000 patches are randomly selected to generate 500 visual words. The codebook is mainly to encode the patches of an image to generate a feature. Then, the features of the training images form the training dictionary. The SPM layer is set to have three different levels of resolution in the proposed LGF, and the nearest 50 atoms in the codebook are used for LLC. As for global features extraction, [4] suggested to set the number of scale to 6, the radius to 3, and the number of neighbor to 10 for MS-CLBP. To generate the Gabor feature images, four directions are used. Then, MS-CLBP operator is applied to these four Gabor feature images.

The proposed LGF strategy is compared with the state-of-the-art LLC + SVM and MS-CLBP + ELM<sup>2</sup> as well as LLC + KCRC and MS-CLBP + KCRC to verify the effectiveness of the proposed fusion approach. For MS-CLBP, the parameters utilized to generate the global features for the training and testing samples are described in [4]. In the experimental setup, all the training and testing features are normalized to [0,1], KCRC is employed with a RBF kernel (optimal parameter  $\sigma = 0.5$ )

<sup>2</sup> ELM: extreme learning machines; the code is available at [http://www3.ntu.edu.sg/home/egbhuang/elm\\_codes.html](http://www3.ntu.edu.sg/home/egbhuang/elm_codes.html).



Fig. 7. Example images from the sports event categories dataset.

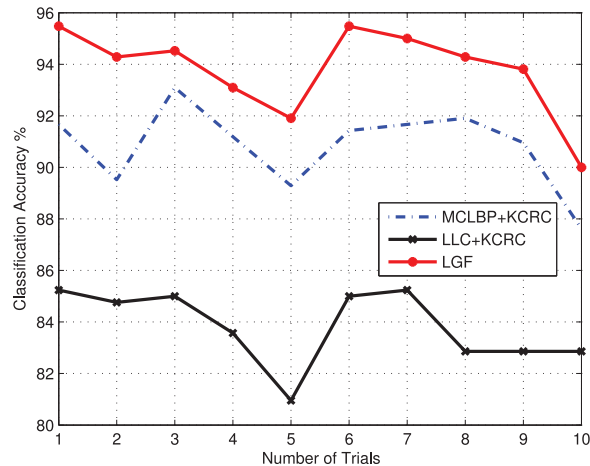


Fig. 8. Classification performance for the UC Merced land-use dataset.

for classification, and the regularization parameter  $\lambda$  is empirically set to  $10^{-4}$ . The weighting parameter  $\mu$  (in Eq. (12)) balancing the global and local feature reconstruction errors varies from 0 to 1 with a step size of 0.1. When  $\mu = 0$ , LGF is equal to MS-CLBP + KCRC; and when  $\mu = 1$ , it reduces to LLC + KCRC.

First of all, we illustrate the classification performance of MS-CLBP + KCRC, LLC + KCRC, and the proposed LGF corresponding to 10 different trials with randomly selected training samples as shown in Fig. 8. It is apparent that the proposed LGF consistently outperforms MS-CLBP + KCRC and LLC + KCRC. In addition, we observe that MS-CLBP + KCRC is superior to LLC + KCRC because this dataset mostly consists of land-use scene images at a large scale and global features are more effective than local features. Cross-validation strategy is employed for tuning the optimal parameter (e.g.,  $\mu$ ) using available training data. In Fig. 9, when  $\mu$  is 0.8, the proposed LGF achieves the highest accuracy, which indicates that the residual residuals from using local features are assigned with a larger weight, leading to better discrimination. This experiment validates the effectiveness of the fusion strategy.

Then, a statistical analysis is conducted by counting the correctly-classified images out of 420 testing images. As listed in Table 1,  $\checkmark$  indicates a correctly-classified sample and  $\times$  indicates a mis-classified sample. The number at each column means the total number of images satisfying three conditions. It exhibits that in this dataset when the LLC + KCRC and MS-CLBP + KCRC results are in ambiguity, the LGF can produce correct results for 65 (columns two and three,  $16 + 50$ ) images out of 72 images (columns two and three, columns six and seven,  $16 + 50 + 7 + 0$ ) and misclassify 7 (columns six and seven,  $7 + 0$ ). When the LLC + KCRC and MS-CLBP + KCRC results are accurate, the proposed LGF will not degrade the results (see the fifth column).



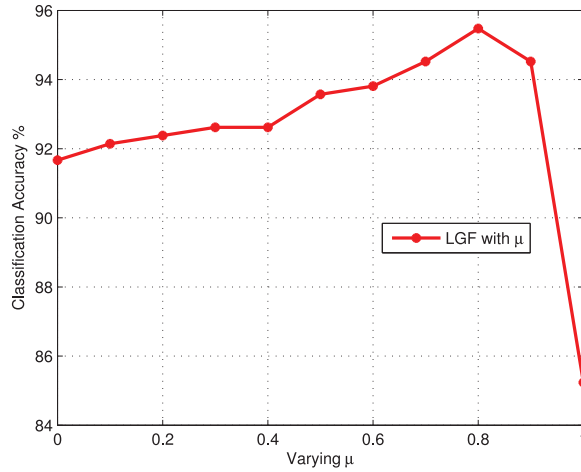


Fig. 9. Different  $\mu$  in LGF for UC Merced land-use dataset.

Table 1

Counting the numbers of the classification accuracy of LLC + KCRC, MS-CLBP + KCRC, and the proposed LGF using the UC Merced land-use dataset.

LLC + KCRC	×	✓	×	✓	✓	✓	×	×
MS-CLBP + KCRC	×	×	✓	✓	✓	×	✓	×
LGF (proposed)	✓	✓	✓	✓	×	×	×	×
	0	16	50	355	0	7	0	12

Table 2

Classification accuracy (%) per class for classifiers using UC Merced land-used dataset.

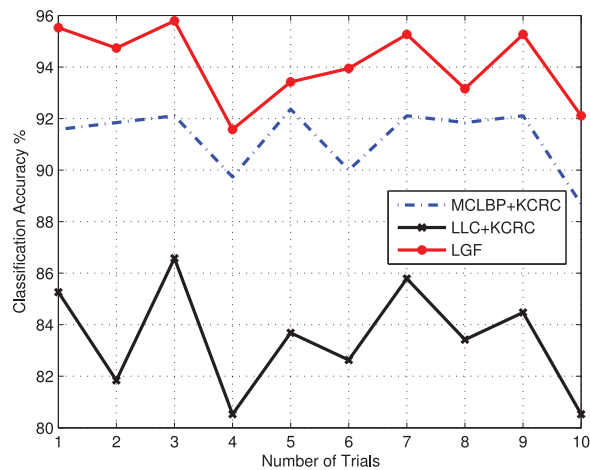
Class	Classification algorithms				
	LLC + SVM	MS-CLBP + ELM	LLC + KCRC	MS-CLBP + KCRC	LGF (proposed)
Agricultural	100	80.00	100	100	100
Airplane	75.00	90.00	80.00	100	100
Baseball-diamond	70.00	85.00	90.00	80.00	100
Beach	100	100	100	100	100
Buildings	65.00	85.00	65.00	85.00	85.00
Chaparral	100	100	100	100	100
Dense-residential	70.00	85.00	75.00	95.00	100
Forest	100	100	100	100	100
Freeway	95.00	75.00	95.00	90.00	100
Golf course	75.00	100	75.00	100	100
Harbor	100	100	100	100	100
Intersection	95.00	95.00	95.00	100	100
Medium-residential	55.00	95.00	75.00	90.00	95.00
Mobile home park	95.00	100	65.00	90.00	95.00
Overpass	80.00	90.00	80.00	70.00	75.00
Parking lot	100	95.00	100	100	100
River	80.00	85.00	75.00	85.00	100
Runway	95.00	95.00	100	100	100
Sparse-residential	65.00	65.00	85.00	95.00	100
Storage tanks	60.00	75.00	55.00	60.00	65.00
Tennis court	70.00	80.00	80.00	85.00	90.00
OA (%)	82.86	89.29	85.24	91.67	<b>95.48</b>
F-measure	0.8293	0.8951	0.8557	0.9189	<b>0.9566</b>

In the first experiment, the accuracy per class from the aforementioned methods and F-measure (calculated by precision and recall [12,29]) is provided in Table 2. Obviously, the proposed LGF achieves better performance than the other methods. It gains almost 13% improvement over the state-of-the-art LLC + SVM method. The KCRC method also outperforms the widely-used SVM classifier. The LGF gains about 4% higher overall accuracy than the global feature representation method, i.e., MS-CLBP + KCRC, and 10% higher than the local feature representation method, i.e., LLC + KCRC. Table 3 displays the

**Table 3**

Confusion matrix for the proposed LGF using UC Merced land-use dataset.

Class	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	Recall(%)
Agricultural (1)	20	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
Airplane (2)	-	20	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
Baseball-diamond (3)	-	-	20	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
Beach (4)	-	-	-	20	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
Buildings (5)	-	-	-	-	17	-	2	-	-	-	-	1	-	-	-	-	-	-	-	-	-	85.00
Chaparral (6)	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
Dense-residential (7)	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
Forest (8)	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	-	-	-	-	100
Freeway (9)	-	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	-	-	-	100
Golf course (10)	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	-	-	100
Harbor (11)	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	-	100
Intersection (12)	-	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	100
Medium-residential (13)	-	-	-	-	-	-	-	-	-	-	-	-	19	1	-	-	-	-	-	-	-	95.00
Mobile home park (14)	-	-	-	-	-	-	-	-	-	-	-	-	1	19	-	-	-	-	-	-	-	95.00
Overpass (15)	-	-	-	-	-	-	1	-	1	-	-	3	-	-	15	-	-	-	-	-	-	75.00
Parking lot (16)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	100
River (17)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	100
Runway (18)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	20	-	-	-	100
Sparse-residential (19)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	20	-	-	100
Storage tanks (20)	-	2	-	-	3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	13	1	65.00
Tennis court (21)	-	1	-	-	-	-	-	-	-	-	-	-	1	-	-	-	-	-	-	-	18	90.00
Precision (%)	100	86.96	100	100	85.00	100	86.96	100	95.24	100	100	83.33	90.48	95.00	100	100	95.24	100	100	100	94.74	

**Fig. 10.** Classification performance in the 19-class satellite scene experiment.

confusion matrix<sup>3</sup> for the proposed LGF, which shows the classification performance (i.e., precision and recall) for individual categories. The major confusion occurs between class 15 (i.e., *overpass*) and class 12 (i.e., *intersection*), or class 20 (i.e., *storage tanks*) and class 5 (i.e., *buildings*).

For the second dataset, the dense SIFT descriptors are extracted for each  $32 \times 32$  patch on a grid of step size 12 since the image size is twice larger than that of the UC-Merced dataset. 500 visual words are generated as a codebook and all the training images as the dictionary. For training, 570 images (30 per class) are chosen, resulting in 1300000 patches, from which 30000 are randomly selected to obtain the codebook. The SPM layer is set to have three different levels of resolution, and 50 nearest atoms in the codebook are used for LLC. As for the MS-CLBP parameters, the number of scale is 6, the radius is 4, and the number of neighbor is 12. Four different directions are employed for generating Gabor feature images,  $\sigma$  is set to 0.9 in KCRC for local feature representation, and 3 for global feature representation after tuning the parameters. The regularization parameter  $\lambda$  is  $10^{-4}$ . Fig. 10 illustrates the results with 10 different trials with random training samples, where the proposed LGF gains the highest accuracy. The global feature can present better than the local feature in this remote sensing dataset because of the large scale of remote sensing images. In this dataset, when  $\mu$  is 0.6, the LGF can obtain the highest fusion accuracy as shown in Fig. 11. To further evaluate the effectiveness of LGF, the number of correctly

<sup>3</sup> The diagonal elements in the matrix indicate the correctly classified number of samples for the classes. The sum of each row indicates the total number of samples for the corresponding class.

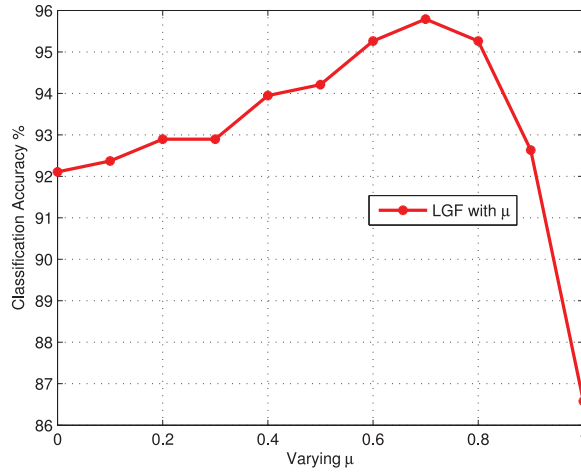


Fig. 11. Different  $\mu$  in LGF for the 19-class satellite scene dataset.

Table 4

Counting the numbers of the classification accuracy of LLC + KCRC, MS-CLBP + KCRC, and the proposed LGF using the 19 class land-use dataset.

LLC + KCRC	×	√	×	√	√	√	×	×
MS-CLBP + KCRC	×	×	√	√	√	×	√	×
LGF (proposed)	√	√	√	√	×	×	×	×
	3	12	31	316	0	1	3	14

Table 5

Classification accuracy (%) per class for classifiers using 19 class land-use dataset.

Class	Classification algorithms				
	LLC + SVM	MSCLBP + ELM	LLC + KCRC	MS-CLBP + KCRC	LGF (proposed)
Airport	90.00	80.00	50.00	80.00	85.00
Beach	80.00	90.00	85.00	90.00	100
Bridge	65.00	70.00	75.00	70.00	75.00
Commercial	100	95.00	80.00	90.00	95.00
Desert	80.00	100	95.00	95.00	95.00
Farmland	75.00	100	80.00	90.00	95.00
Forest	95.00	100	100	95.00	100
Industrial	45.00	80.00	95.00	95.00	100
Meadow	95.00	100	60.00	75.00	85.00
Mountain	95.00	100	95.00	100	100
Park	100	100	100	100	100
Parking	100	90.00	100	95.00	100
Pond	100	95.00	100	95.00	100
Port	70.00	95.00	100	100	100
Residential	75.00	95.00	75.00	95.00	95.00
River	75.00	100	90.00	90.00	100
Viaduct	80.00	90.00	80.00	100	100
Football field	100	90.00	85.00	95.00	85.00
Railway station	90.00	80.00	100	100	100
OA (%)	84.73	92.10	86.58	92.11	<b>95.26</b>
F-measure	0.8418	0.9227	0.8671	0.9246	<b>0.9591</b>

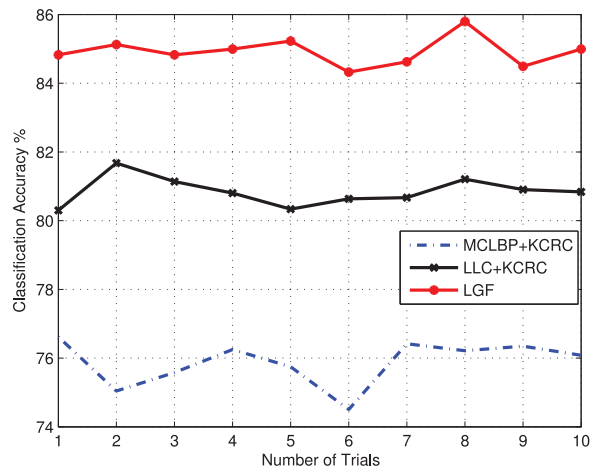
classified images (out of 380 testing images) is counted. As listed in Table 4, when both of the LLC+KCRC and MS-CLBP + KCRC produce wrong labels, the LGF can correct 3 images from 17 as shown in the first column. Table 5 lists the accuracy of each class for different methods, where the proposed LGF achieves the highest accuracy and outperforms the second best result from the MS-CLBP + ELM by 3%. Table 6 further shows the confusion matrix for the proposed LGF.

To evaluate the proposed method on the nature scene image datasets in the third experiment, all the images size are resized to 300 × 300 (some of them may be smaller) to reduce computational complexity. The dense SIFT descriptors are similarly extracted. The size of codebook is set to 1024 because a large visual word codebook is necessary for this dataset consisting of both indoor and outdoor images. For training, 1500 images (100 per class) are used, generating 2,880,000

**Table 6**

Confusion matrix for the proposed LGF using 19 class land-use dataset.

Class	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	Recall (%)
Airport (1)	17	-	-	1	-	-	1	-	-	-	-	-	-	-	-	-	-	-	1	85.00
Beach (2)	-	20	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
Bridge (3)	-	-	15	-	-	-	-	-	-	1	-	-	-	1	1	-	-	2	-	75.00
Commercial (4)	-	-	-	19	-	-	-	-	-	-	-	-	-	-	-	-	1	-	-	95.00
Desert (5)	-	-	-	-	19	-	-	-	-	-	-	-	-	-	1	-	-	-	-	95.00
Farmland (6)	-	-	-	-	-	19	-	-	-	-	-	-	-	-	-	1	-	-	-	95.00
Forest (7)	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	-	-	-	100
Industrial (8)	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	-	-	100
Meadow (9)	-	-	-	-	-	-	-	-	17	-	-	-	-	-	-	-	2	-	1	85.00
Mountain (10)	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	-	100
Park (11)	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	-	100
Parking (12)	-	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	-	-	100
Pond (13)	-	-	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	-	100
Port (14)	-	-	-	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	-	100
Residential (15)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	20	-	-	-	-	100
River (16)	-	-	-	-	-	-	-	-	1	-	-	-	-	-	-	19	-	-	-	95.00
Viaduct (17)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	20	-	-	100
Football field (18)	-	-	-	-	-	-	-	1	-	-	-	-	-	-	-	-	-	19	-	95.00
Railway station (19)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	20	100
Precision (%)	100	100	100	95.00	100	100	95.24	95.24	94.44	95.24	100	100	100	95.24	90.91	95.00	86.96	90.48	90.91	

**Fig. 12.** Classification performance in the 15-scene indoor and outdoor data experiment.**Table 7**

Counting the numbers of the classification accuracy of LLC + KCRC, MS-CLBP + KCRC, and the proposed LGF using the 15-scene categories.

LLC + KCRC	x	√	x	√	√	√	x	x
MS-CLBP + KCRC	x	x	√	√	√	x	√	x
LGF (proposed)	√	√	√	√	x	x	x	x
	27	328	203	2003	0	93	69	262

patches, from which 30,000 patches are randomly selected to construct the codebook by  $k$ -means clustering. The SPM layer is set to have three different levels of resolution, and the nearest 50 atoms in the codebook are calculated for LLC. For the global feature parameter, the scale is set to 6 for MS-CLBP, the radius is 3, and the neighboring point is 10 for the CLBP operator. Four different directions are used for generating Gabor feature images. In the proposed LGF, the RBF kernel parameter of the KCRC is  $\sigma = 1$  for local features and  $\sigma = 0.5$  for global features according to our empirical study, and the regularized parameter  $\lambda$  is  $10^{-4}$ . Experimental results with 10 different trials with random training samples are shown in Fig. 12, where it is apparent that the LGF outperforms the MS-CLBP + KCRC and LLC + KCRC. From the results, the performance of LLC+KCRC is better than MS-CLBP + KCRC, which is mainly due to the fact that these nature scene images present local details. It is reasonable that the local feature based methods perform better than the global feature based methods. Fig. 13 illustrates that when  $\mu$  is 0.8, the proposed method achieves the highest accuracy. Table 7 demonstrates that when both the results from the LLC + KCRC and MS-CLBP + KCRC are wrong, the LGF can correct 27 images out of

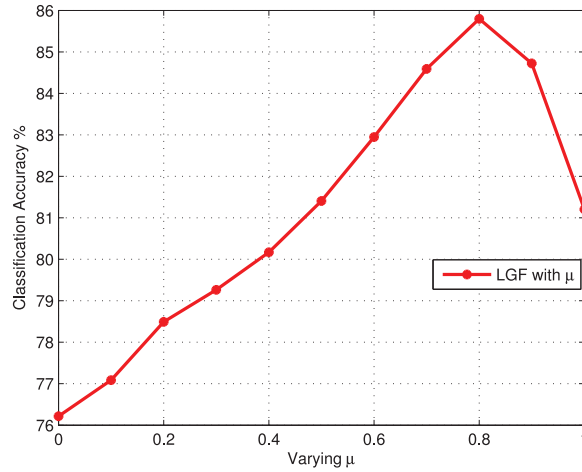


Fig. 13. Different  $\mu$  in LGF for 15-scene categories.

Table 8  
Classification accuracy (%) per class for classifiers using 15-scene categories.

Class	Classification algorithms				
	LLC + SVM	MSCLBP + ELM	LLC + KCRC	MS-CLBP + KCRC	LGF (proposed)
CAL-suburb	99.28	99.28	98.58	96.45	100
MIT-coast	90.77	77.31	87.31	79.62	88.85
MIT-forest	92.98	85.53	96.49	89.91	95.61
MIT-highway	86.25	80.00	88.13	85.00	91.25
MIT-insidecity	82.21	68.75	81.25	66.35	83.65
MIT-mountain	90.15	79.20	86.86	81.75	90.15
MIT-opencountry	68.06	66.77	72.58	69.68	74.52
MIT-street	90.10	76.04	91.15	85.42	91.67
MIT-tallbuilding	90.63	81.25	89.84	70.70	92.97
PARoffice	95.65	95.65	95.65	95.65	99.13
Bedroom	68.97	52.59	68.97	56.03	68.97
Industrial	51.66	58.77	46.92	69.67	72.99
Kitchen	72.73	59.09	72.73	67.27	82.73
Livingroom	69.31	44.44	69.31	60.85	71.96
Store	73.49	72.56	74.42	73.02	85.58
OA (%)	81.34	73.20	81.21	76.21	<b>85.80</b>
F-measure	0.8117	0.7230	0.8129	0.7579	<b>0.8572</b>

288. When one result from the LLC + KCRC and MS-CLBP + KCRC is wrong, the LGF can produce 531 correct labels and only remaining 162 images are misclassified. Eventually, 2561 out of 2985 images are assigned with correct labels. The accuracy of each class is listed in Table 8, where the proposed LGF method yields 4% higher accuracy than the second best result. Table 9 further displays the confusion matrix for the proposed LGF.

The last experiment is carried out on the sports event categories dataset with high resolution, and all the images are resized to  $300 \times 300$  for saving computational efficiency. The dense SIFT descriptors are extracted and the codebook is set to 1024. For training purpose, 560 images (70 per class) are used, and 30,000 patches are randomly selected to generate the codebook. The SPM layer is set to have three different levels of resolution, and the nearest 50 atoms in the codebook are calculated for LLC. For the global feature parameters, the same parameters are used for the other datasets, i.e., scale = 6, radius = 3, and the number of neighbor point = 10 for CLBP. Four different direction Gabor feature images are also employed by CLBP. For the proposed approach, the RBF kernel parameter of the KCRC is  $\sigma = 1$ , and the regularized parameter  $\lambda$  is  $10^{-4}$ . Experimental results with 10 trails with random training samples are illustrated in Fig. 14, where the LGF gains a large margin of improvements over the other algorithms. In this dataset, the LLC + KCRC and MS-CLBP + KCRC produce similar accuracies. This is mainly because almost all the images are with a large scale background and some detailed human actions or objects. As a consequence, both the local and global features can well represent the image features. Fig. 15 shows that when  $\mu$  is 0.7, the LGF yields the highest fusion accuracy. According to Table 10, when both the LLC + KCRC and MSCLBP + KCRC assign incorrect labels, the LGF can correct 5 out of 78 images as indicated in the first column. 136 images are assigned with correct labels and 44 are left with incorrect labels when only one of the LLC + KCRC and MS-CLBP + KCRC gives the right label. Table 11 provides the accuracy of each class from different methods, where the LGR generally gains 4% improvement in accuracy. Table 12 further displays the confusion matrix for the proposed LGF.

**Table 9**  
Confusion matrix for the proposed LGF using 15-scene categories.

Class	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	Recall (%)
CALsuburb (1)	141	-	-	-	-	-	-	-	-	-	-	-	-	-	-	100
MIT-coast (2)	-	231	1	7	-	6	15	-	-	-	-	-	-	-	-	88.85
MIT-forest (3)	-	-	218	-	-	8	2	-	-	-	-	-	-	-	-	95.61
MIT-highway (4)	-	2	-	146	2	3	3	1	-	-	-	1	1	-	1	91.25
MIT-insidicity (5)	2	-	-	2	174	-	-	13	12	-	-	1	2	1	1	83.65
MIT-mountain (6)	-	6	7	-	1	247	11	-	1	-	-	-	-	-	1	90.15
MIT-opencountry (7)	-	39	13	8	-	17	231	1	-	-	-	1	-	-	-	74.52
MIT-street (8)	-	-	-	3	5	2	-	176	4	-	-	1	-	1	-	91.67
MIT-tallbuilding (9)	-	1	3	-	6	1	-	1	238	-	-	2	-	2	2	92.97
PARoffice (10)	-	-	-	-	-	-	-	-	-	114	-	-	1	-	-	99.13
Bedroom (11)	2	1	-	-	-	1	-	1	-	7	80	-	4	19	1	68.97
Industrial (12)	2	1	1	-	7	1	3	1	8	1	1	154	3	4	24	72.99
Kitchen (13)	-	-	-	-	2	1	-	-	-	1	3	-	91	8	4	82.73
Livingroom (14)	-	-	-	-	3	1	-	3	-	5	26	1	9	136	5	71.96
Store (15)	-	-	4	-	6	2	-	-	3	2	2	2	6	4	184	85.58
Precision (%)	95.92	82.21	88.26	87.95	84.47	85.17	87.17	89.34	89.47	87.69	71.43	94.48	77.78	77.71	82.51	

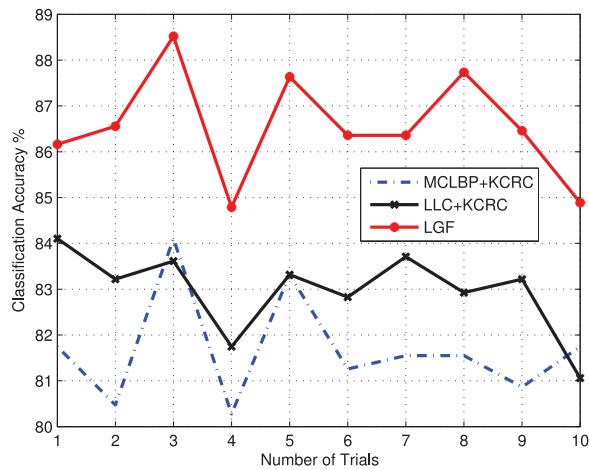


Fig. 14. Classification performance in the sports event experiment.

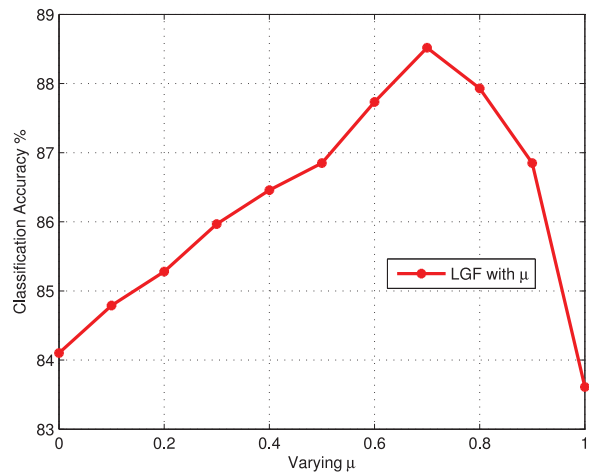


Fig. 15. Different μ in LGF for sports event categories.

**Table 10**

Counting the numbers of the classification accuracy of LLC + KCRC, MS-CLBP + KCRC, and the proposed LGF using the sports event categories.

LLC + KCRC	×	✓	×	✓	✓	✓	×	×
MS-CLBP + KCRC	×	×	✓	✓	✓	×	✓	×
LGF (proposed)	✓	✓	✓	✓	×	×	×	×
	5	58	71	768	0	26	18	73

**Table 11**

Classification accuracy (%) per class for classifiers using sports event categories.

Class	Classification algorithms				
	LLC + SVM	MSCLBP + ELM	LLC + KCRC	MS-CLBP + KCRC	LGF (proposed)
RockClimbing	91.13	99.28	96.77	93.55	95.97
Badminton	90.00	77.31	95.38	94.62	96.92
Bocce	55.22	85.53	56.72	68.66	68.66
Croquet	75.30	80.00	69.28	82.53	84.34
Polo	85.71	68.75	84.82	75.89	83.93
Rowing	87.22	79.20	87.22	80.56	88.89
Sailing	91.67	66.77	91.67	90.00	95.83
Snowboarding	80.00	76.04	77.50	80.83	85.00
OA (%)	83.51	73.20	83.61	84.10	<b>88.52</b>
F-measure	0.8189	0.7129	0.8215	0.8299	<b>0.8713</b>

**Table 12**

Confusion matrix for the proposed LGF using sports event categories.

Class	1	2	3	4	5	6	7	8	Recall (%)
RockClimbing (1)	119	–	–	–	1	–	–	4	95.97
Badminton (2)	–	126	2	2	–	–	–	–	96.92
Bocce (3)	4	3	46	6	2	1	2	3	68.66
Croquet (4)	5	–	16	140	3	–	1	1	84.34
Polo (5)	2	2	2	2	94	2	2	6	83.93
Rowing (6)	2	–	4	1	3	160	5	5	88.89
Sailing (7)	–	–	–	3	–	1	115	1	95.83
Snowboarding (8)	6	3	1	1	3	2	2	102	85.00
Precision (%)	86.23	94.03	64.79	90.32	88.68	96.39	90.55	83.61	

**Table 13**

Comparison of classification accuracy on the 21-class land-use dataset.

Algorithms	Accuracy (%)
BOVW [54]	76.8
SPM [54]	75.3
BOVW + Spatial co-occurrence kernel [54]	77.7
Color Gabor [54]	80.5
Color histogram [54]	81.2
Structural texture similarity [39]	86.0
Wavelet BOVW [65]	87.4
Unsupervised feature learning [7]	81.7
Saliency-guided feature learning [60]	82.7
Concentric circle-structured BOVW [66]	86.6
Multifeature concatenation [43]	89.5
Pyramid-of-spatial-relations [6]	89.1
MS-CLBP + ELM [4]	89.29
LLC + SVM	82.86
MS-CLBP + KCRC	91.67
LLC + KCRC	85.24
Proposed LGF	<b>95.48</b>

The proposed LGF is also compared with several state-of-the-art algorithms in the literature. As reported in Tables 13–15, the LGF can improve the accuracy by 4% even based on the original features. In addition, the proposed fusion strategy is also suitable for combining different features.

Finally, computing time on the environment of Linux with a Core I7 CPU and 8G memory is listed in Table 17. It indicates that local feature extraction has more complexity than global feature extraction because of more float computing in patches and coding with BoVW, while the global features computing in entire images and only mapping to bins. According to our experiments, the local feature extracting time needs about 0.7 s for images not large than 300 pixels and almost 2 s per

**Table 14**

Comparison of classification accuracy on the 19-class satellite scene dataset.

Algorithms	Accuracy (%)
Bag of colors [43]	70.6
Tree of c-shapes [43]	80.4
Bag of SIFT [43]	85.5
Multifeature concatenation [43]	90.8
LTP-HF [44]	77.6
SIFT + LTP-HF + color histogram [44]	93.6
MS-CLBP + KML [4]	92.71
LLC + SVM	82.75
MS-CLBP + KCRC	92.11
LLC + KCRC	86.58
Proposed LGF	<b>95.26</b>

**Table 15**

Comparison of classification accuracy on the 15-scene categories.

Algorithms	Accuracy (%)
KSPM [16]	81.40
KC [47]	76.67
KSPM [52]	76.73
LSPM [52]	65.32
ScSPM [52]	80.28
MS-CLBP + KML [4]	70.13
LLC + SVM	81.34
MS-CLBP + KCRC	76.21
LLC + KCRC	81.21
Proposed LGF	<b>85.80</b>

**Table 16**

Comparison of classification accuracy on the sports event categories.

Algorithms	Accuracy (%)
Full integrative model [19]	73.4
OB [21]	76.3
SIFT + SC [52]	82.7
HMP [2]	85.7
MS-CLBP + KML [4]	73.20
LLC + SVM	83.51
MS-CLBP + KCRC	84.10
LLC + KCRC	83.61
Proposed LGF	<b>88.52</b>

**Table 17**

The cost complexity (in seconds) of the proposed method for all the experimental data.

	UC Merced land-use	19-class satellite	15-scene indoor and outdoor	Sport event categories
LLC + KCRC	302.4	760.1	2029.8	764.2
MS-CLBP + KCRC	71.4	304.1	477.6	173.7
LGF (proposed)	365.6	1054.3	2512.6	948.7

image with  $600 \times 600$  pixels. The global features need only 0.2 s for images not large than 300 pixels and 0.8 s for an image of size  $600 \times 600$ . The proposed LGF needs to compute both local and global features, and represent and reconstruct in classification stage, resulting in higher time complexity. However, it is also noticed that the proposed algorithm is feasible to operate by parallel computation; in doing so, computing complexity will have little side-effect.

## 5. Conclusion

In this paper, we proposed an effective feature extraction and representation-based fusion strategy for scene classification. The proposed method combines local feature representation based on BoVM and SPM and global feature representation based on MS-CLBP. It can well explore the complementary nature of local and global features, leading to enhanced feature discriminative power. A weighted sum of reconstruction residuals from using local and global features in collaborative representation were then developed for image classification. The proposed approach was thoroughly evaluated on four different



datasets including remote sensing images, indoor and outdoor nature images, and human sports action images. The experimental results demonstrated the superiority of the proposed method over several state-of-the-art scene classification methods.

## Acknowledgment

This work was supported by the National Natural Science Foundation of China under Grants no. NSFC-61571033, 61302164, and partly by the Fundamental Research Funds for the Central Universities under Grants no. BUCTRC201401, YS1404, XK1521.

## References

- [1] T. Ahonen, A. Hadid, M. Pietikinen, Face description with local binary patterns: Application to face recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 28 (12) (2006) 2037–2041.
- [2] L. Bo, X. Ren, D. Fox, Hierarchical matching pursuit for image classification: Architecture and fast algorithms, in: *Advances in Neural Information Processing Systems*, 2011.
- [3] S. Chaib, Y. Gu, H. Yao, An informative feature selection method based on sparse PCA for VHR scene classification, *IEEE Geosci. Remote Sens. Lett.* 13 (2) (2016) 147–151.
- [4] C. Chen, B. Zhang, H. Su, W. Li, L. Wang, Land-use scene classification using multi-scale completed local binary patterns, *Signal Image Video Process.* (2015), doi:10.1007/s11760-015-0804-2.
- [5] C. Chen, L. Zhou, J. Guo, W. Li, H. Su, F. Guo, Gabor-filtering-based completed local binary patterns for land-use scene classification, in: *Proceedings of the First IEEE International Conference on Multimedia Big Data and Hyperspectral Imaging Workshop*, 2015.
- [6] S. Chen, Y. Tian, Pyramid of spatial relations for scene-level land use classification, *IEEE Trans. Geosci. Remote Sens.* 53 (4) (2015) 1947–1957.
- [7] A.M. Cheryadat, Unsupervised feature learning for aerial scene classification, *IEEE Trans. Geosci. Remote Sens.* 52 (1) (2014) 439–451.
- [8] D. Das, Scene classification using pyramid histogram of multi-scale block local binary pattern, *Int. J. Comput. Sci. Appl. (IJCSA)* 4 (4) (2014) 15–25.
- [9] D. Lowe, Distinctive image features from scale-invariant keypoints, *Int. J. Comput. Vis.* 60 (2004) 91–110.
- [10] L. Duan, J. Lin, Z. Wang, T. Huang, W. Gao, Weighted component hashing of binary aggregated descriptors for fast visual search, *IEEE Trans. Multimed.* 17 (6) (2015) 828–842.
- [11] K. Grauman, T. Darrell, Pyramid match kernels: Discriminative classification with sets of image features, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2005.
- [12] J.-M. Guo, C.-H. Hsia, Y.-F. Liu, M.-H. Shih, C.-H. Chang, J.-Y. Wu, Fast background subtraction based on a multiplayer codebook model for moving object detection, *IEEE Trans. Circuits Syst. Video Technol.* 23 (10) (2013) 1809–1821.
- [13] Z. Guo, L. Zhang, D. Zhang, A completed modeling of local binary pattern operator for texture classification, *IEEE Trans. Image Process.* 19 (6) (2010) 1657–1663.
- [14] F. Hu, G. Xia, Z. Wang, X. Huang, L. Zhang, H. Sun, Unsupervised feature learning via spectral clustering of multidimensional patches for remotely sensed scene classification, *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 8 (5) (2015) 2015–2030.
- [15] Y. Jiang, J. Wang, X. Xue, S.-F. Chang, Query adaptive image search with hash codes, *IEEE Trans. Multimed.* 15 (2) (2013) 442–453.
- [16] S. Lazebnik, C. Schmid, J. Ponce, beyond bags of features: Spatial pyramid matching for recognizing natural scene categories, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2006, pp. 2169–2178.
- [17] B. Li, D. Ming, W. Yan, X. Sun, T. Tian, J. Tian, Image matching based on two-column histogram hashing and improved RANSAC, *IEEE Geosci. Remote Sens. Lett.* 11 (8) (2014) 1433–1437.
- [18] F.-F. Li, P. Perona, A Bayesian hierarchical model for learning natural scene categories, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2005, pp. 524–531.
- [19] L. Li, F.-F. Li, What, where and who? classifying events by scene and object recognition, in: *Proceedings of the Eleventh IEEE International Conference on Computer Vision*, 2007, pp. 1–8.
- [20] L. Li, H. Su, Y. Lim, F.-F. Li, Objects as attributes for scene classification, in: *Proceedings of the European Conference Computer Vision*, 2010, pp. 1–13.
- [21] L. Li, H. Su, E. Xing, F.-F. Li, Object bank: A high-level image representation for scene classification and semantic feature sparsification, in: *Proceedings of the Neural Information Processing Systems*, Vancouver, Canada, 2008.
- [22] W. Li, C. Chen, H. Su, Q. Du, Local binary patterns and extreme learning machine for hyperspectral imagery classification, *IEEE Trans. Geosci. Remote Sens.* 53 (7) (2015) 3681–3693.
- [23] W. Li, Q. Du, Collaborative representation for hyperspectral anomaly detection, *IEEE Trans. Geosci. Remote Sens.* 53 (3) (2015) 1463–1474.
- [24] Y. Li, C. Tao, Y. Tan, K. Shang, J. Tian, Unsupervised multilayer feature learning for satellite image scene classification, *IEEE Geosci. Remote Sens. Lett.* 13 (2) (2016) 157–161.
- [25] J. Liu, S. Mubarak, Scene modeling using co-clustering, in: *Proceedings of the Eleventh IEEE International Conference on Computer Vision*, 2007, pp. 1–7.
- [26] J. Liu, Z. Wu, Z. Wei, L. Xiao, L. Sun, Spatial-spectral kernel sparse representation for hyperspectral image classification, *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 6 (6) (2013) 2462–2471.
- [27] D. Lowe, Object recognition from local scale-invariant features, *Int. J. Comput. Vis.* 2 (1999) 1150–1157.
- [28] D.G. Lowe, Local feature view clustering for 3d object recognition, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2001, pp. 682–688.
- [29] S. Marina, G. Lapalme, A systematic analysis of performance measure for classification tasks, *Inf. Process. Manag.* 45 (4) (2009) 427–437.
- [30] K. Mikolajczyk, C. Schmid, Scale and affine invariant interest point detectors, *Int. J. Comput. Vis.* 60 (2004) 63–86.
- [31] K. Mikolajczyk, C. Schmid, A performance evaluation of local descriptors, *IEEE Trans. Pattern Anal. Mach. Intell.* 27 (2005) 1615–1630.
- [32] K.-R. Muller, S. Mika, G. Ratsch, K. Tsuda, B. Scholkopf, An introduction to kernel-based learning algorithms, *IEEE Trans. Neural Netw.* 12 (2) (2001) 181–201.
- [33] T. Ojala, M. Pietikainen, T.T. Maenpaa, Multiresolution gray-scale and rotation invariant texture classification with local binary patterns, *IEEE Trans. Pattern Anal. Mach. Intell.* 24 (7) (2002) 971–987.
- [34] A. Oliva, A. Torralba, Modeling the shape of the scene: a holistic representation of the spatial envelope, *Int. J. Comput. Vis.* 42 (3) (2001) 145–175.
- [35] B. Peng, W. Li, X. Xie, Q. Du, K. Liu, Weighted fusion-based representation classifiers for hyperspectral imagery, *Remote Sens.* 7 (11) (2015) 14806–14826.
- [36] K. Qi, H. Wu, C. Shen, J. Gong, Land-use scene classification in high-resolution remote sensing images using improved correlations, *IEEE Geosci. Remote Sens. Lett.* 12 (12) (2015) 2403–2407.
- [37] J. Qin, N.H.C. Yung, Scene categorization via contextual visual words, *Pattern Recognit.* 43 (2010) 1874–1888.
- [38] A. Quattoni, A. Torralba, Recognizing indoor scenes, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2009, pp. 413–420.
- [39] V. Risojevic, Z. Babic, Aerial image classification using structural texture similarity, in: *Proceedings of the IEEE International Symposium on Signal Processing and Information Technology*, 2011, pp. 190–195.

- [40] V. Risojevic, Z. Babic, Fusion of global and local descriptors for remote sensing image classification, *IEEE Geosci. Remote Sens. Lett.* 10 (4) (2013) 836–840.
- [41] N. Serrano, A. Savakis, J. Luo, Improved scene classification using efficient low-level features and semantic cues, *Pattern Recognit.* 37 (2004) 1773–1784.
- [42] T. Serre, M. Kouh, C. Cadieu, U. Knoblich, G. Kreiman, T. Poggio, A theory of object recognition: Computations and circuits in the feedforward path of the ventral stream in primate visual cortex, in: *Computer Science and Artificial Intelligence Laboratory Technical Report*, 2005.
- [43] W. Shao, W. Yang, G. Xia, G. Liu, A hierarchical scheme of multiple feature fusion for high-resolution satellite scene categorization, in: *Proceedings of the Ninth International Conference on Computer Vision Systems*, St. Petersburg, Russia, 2013, pp. 324–333.
- [44] G. Sheng, W. Yang, T. Xu, H. Sun, High-resolution satellite scene classification using a sparse coding based multiple feature combination, *Int. J. Remote Sens.* 33 (8) (2011) 2395–2412.
- [45] A. Torralba, Contextual priming for object detection, *Int. J. Comput. Vis.* 53 (2) (2003) 169–191.
- [46] A. Torralba, A. Oliva, Statistics of natural image categories, *Netw. Comput. Neural Syst.* 14 (2003) 391–412.
- [47] J.C. van Gemert, J.-M. Geusebroek, C.J. Veenman, A.M. Smeulders, Kernel codebooks for scene categorization, in: *Proceedings of the Tenth European Conference on Computer Vision ECCV*, 2008.
- [48] C. Wallraven, B. Caputo, A. Graf, Recognition with local features: the kernel recipe, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 1, 2003, pp. 257–264.
- [49] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, T. Gong, Locality-constrained linear coding for image classification, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2010, pp. 3360–3367.
- [50] J. Willamowski, D. Arregui, G. Csurka, C.R. Dance, L. Fan, Categorizing nine visual classes using local appearance descriptors, in: *Proceedings of the ICPR Workshop on Learning for Adaptable Visual Systems*, 2004.
- [51] J. Yang, A.F. Frangi, J.Y. Yang, D. Zhang, KPCA plus LDA: a complete kernel fisher discriminant framework for feature extraction and recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 27 (2) (2005) 230–244.
- [52] J. Yang, K. Yu, Y. Gong, T. Huang, Linear spatial pyramid matching using sparse coding for image classification, in: *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2009.
- [53] W. Yang, Z. Wang, J. Yin, C. Sun, K. Ricanek, Image classification using kernel collaborative representation with regularized least square, *Appl. Math. Comput.* 222 (1) (2013) 13–28.
- [54] Y. Yang, S. Newsam, Bag-of-visual-words and spatial extensions for land-use classification, in: *Proceedings of the Eighteenth ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, 2010, pp. 270–279.
- [55] Y. Yang, S. Newsam, Spatial pyramid co-occurrence for image classification, in: *Proceedings of the International Conference on Computer Vision*, 2011, pp. 1465–1472.
- [56] J. Yu, R. Hong, M. Wang, J. You, Image clustering based on sparse patch alignment framework, *Pattern Recognit.* 47 (11) (2014) 3512–3519.
- [57] J. Yu, Y. Rui, Y. Tang, D. Tao, High-order distance-based multiview stochastic learning in image classification, *IEEE Trans. Cybern.* 44 (12) (2014) 2431–2442.
- [58] J. Yu, D. Tao, J. Li, J. Cheng, Semantic preserving distance metric learning and applications, *Inf. Sci.* 281 (2014) 671–686.
- [59] F. Zhang, B. Du, L. Zhang, Saliency-guided unsupervised feature learning for scene classification, *IEEE Trans. Geosci. Remote Sens.* 53 (4) (2015) 2175–2184.
- [60] F. Zhang, B. Du, L. Zhang, Saliency-guided unsupervised feature learning for scene classification, *IEEE Trans. Geosci. Remote Sens.* 53 (4) (2015) 2175–2184.
- [61] J. Zhang, M. Marszalek, S. Lazebnik, C. Schmid, Local Features and Kernels for Classification of Texture and Object Categories: An In-Depth Study, INRIA Rhone-Alpes, 2005. Technical Report RR-5737.
- [62] L. Zhang, M. Yang, X. Feng, Sparse representation or collaborative representation: which helps face recognition, in: *Proceedings of the IEEE Conference on Computer Vision (ICCV2011)*, 2011.
- [63] L. Zhang, X. Zhen, L. Shao, Learning object-to-class kernels for scene classification, *IEEE Trans. Image Process.* 23 (8) (2014) 3241–3253.
- [64] W. Zhang, S. Shan, W. Gao, X. Chen, H. Zhang, Local Gabor binary pattern histogram sequence (LGBPHS): a novel non-statistical model for face representation and recognition, in: *Proceedings of the Tenth IEEE International Conference on ICCV*, 2005, pp. 786–791.
- [65] L. Zhao, L.Huo, A 2-d wavelet decomposition-based bag-of-visual-words model for landuse scene classification, *Int. J. Remote Sens.* 35 (6) (2014) 2296–2310.
- [66] L. Zhao, P. Tang, L.Huo, Land-use scene classification using a concentric circle-structured multiscale bag-of-visual-words model, *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.* 7 (12) (2014) 4620–4631.
- [67] Y. Zhong, Q. Zhu, L. Zhang, Scene classification based on the multifeature fusion probabilistic topic model for high spatial resolution remote sensing imagery, *IEEE Trans. Geosci. Remote Sens.* 53 (11) (2015) 6207–6222.
- [68] L. Zhou, D. Hu, Z. Zhou, Scene recognition combining structural and textural features, in: *Science China Information Sciences*, 2011.
- [69] L. Zhou, Z. Zhou, D. Hu, Scene classification using a multi-resolution bag-of-features model, *Pattern Recognit.* 46 (1) (2013) 424–433.
- [70] L. Zhou, Z. Zhou, D. Hu, Scene classification using a multi-resolution low-level feature combination, *Neurocomput.* 122 (25) (2013) 284–297.